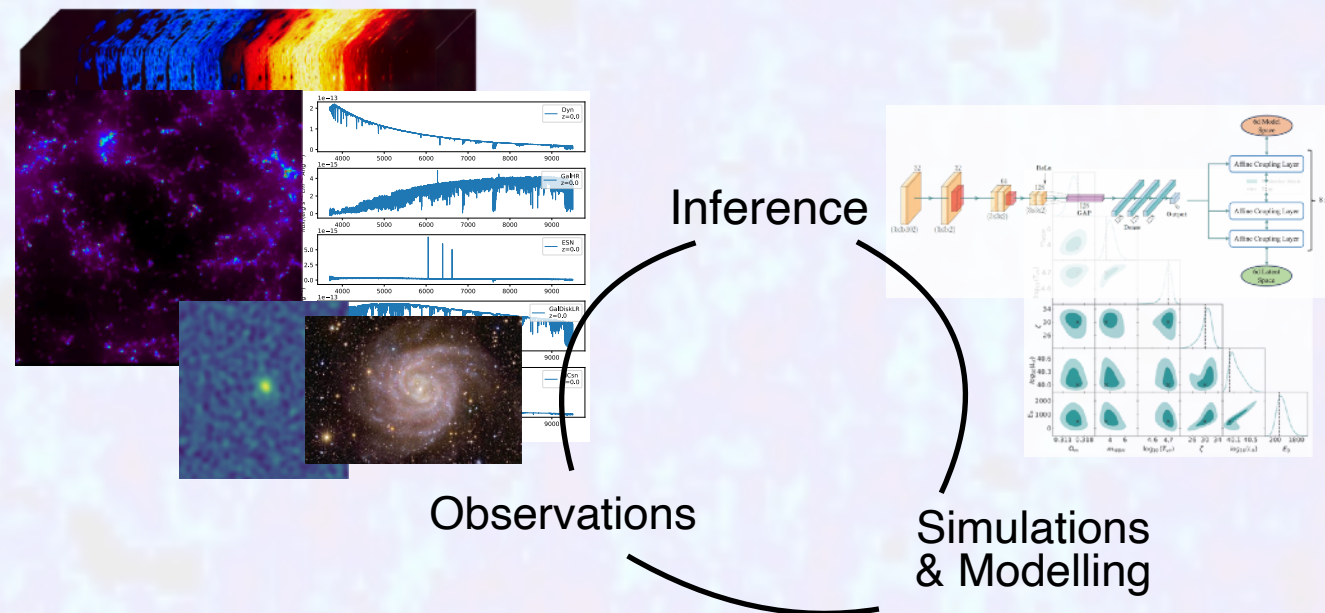# Machine Learning for Astrophysics & Cosmology



Caroline Heneka, group leader, ITP Heidelberg

'Computer Vision Astrophysics and Cosmology'

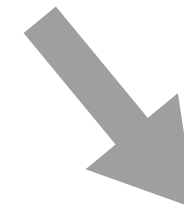Physics in the AI era, Pisa, September 25th 2024

# Cosmology & ML group @ ITP

**Our goal**
Learn about cosmology, large-scale structure, the high-redshift Universe (Reionization) and develop the suitable modern ML toolkit.

**About myself:**
B.Sc. and M.Sc. in Physics (Heidelberg)
PhD in Physics 2017 (Copenhagen)
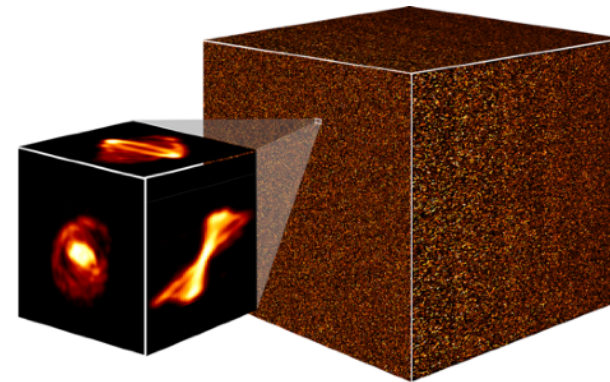Postdoc: SNS Pisa, UHH Hamburg, + DLR
Since 10/22 Group Leader

Our research:
- Computational astrophysics / cosmology
- Intensity mapping
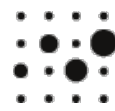- Large (radio) surveys, SKA & LOFAR

Specifically the modern ML toolkit for cosmology and large-scale surveys:
- Emulation, generation
- Inference
- Classification, anomaly detection
- Computer Vision tasks in astronomy

@SKAO

STRUCTURES CLUSTER OF EXCELLENCE

**Daimler** und **Benz** Stiftung

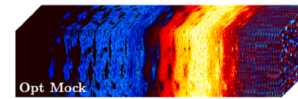**VolkswagenStiftung** FREIGEIST FELLOWSHIP DER VOLKSWAGENSTIFTUNG

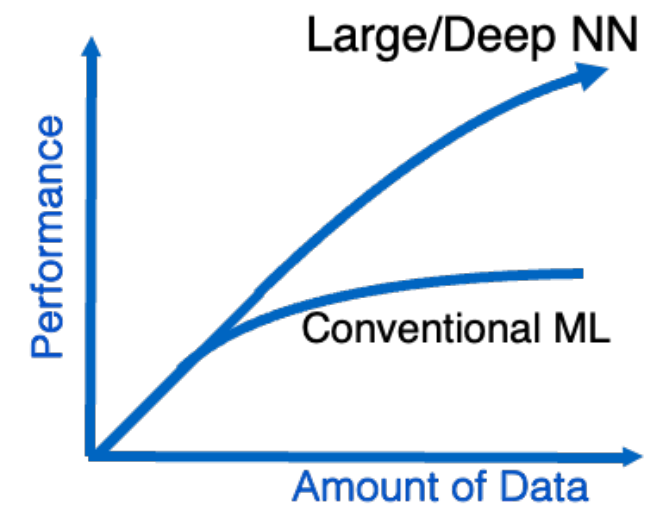# How will Astrophysics and Cosmology advance in coming years?

Understanding



Survey Data

Data Science
'ML + AI'
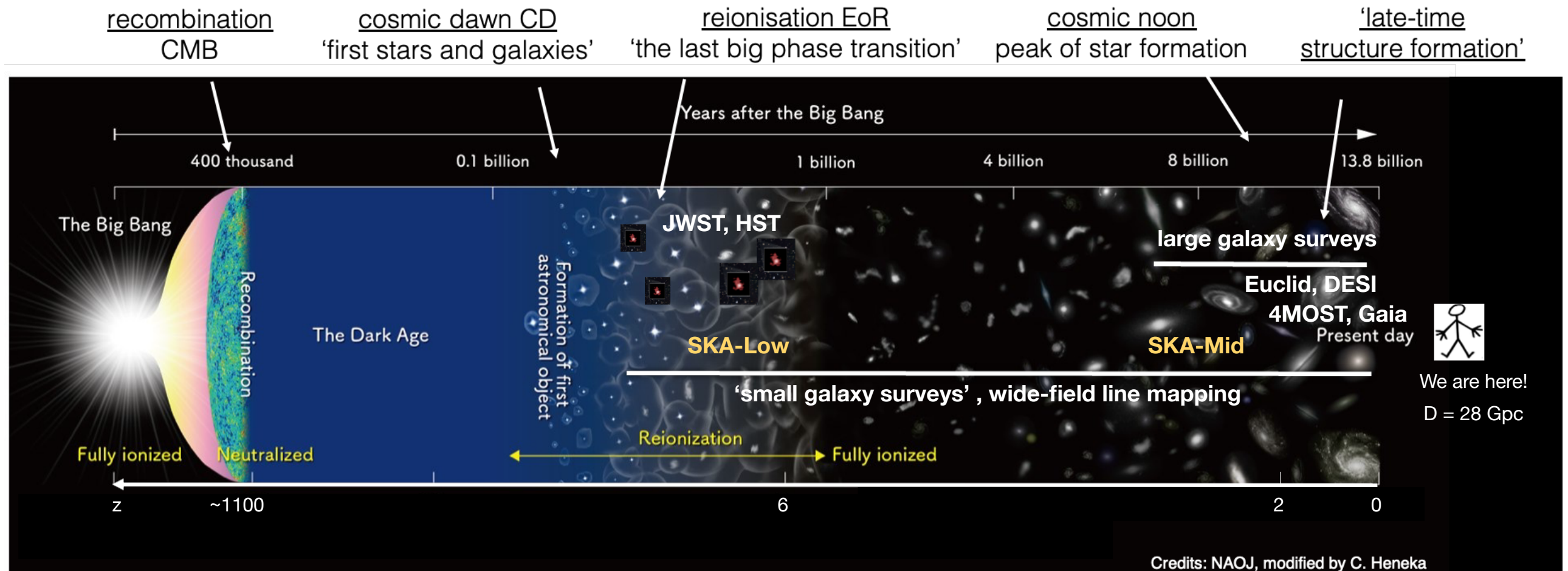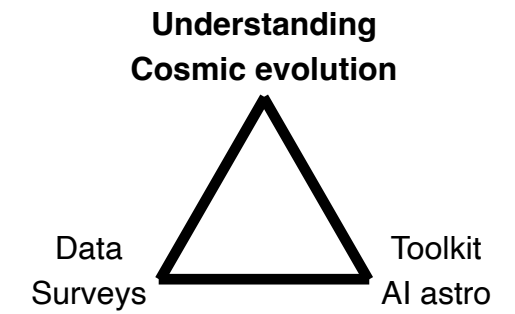
ELT, Credit: ESO



Credit: NASA,

MeerKAT, Credit: SRAO

Credit: ESA

Credit: SKAO

Efficient data reduction
Automation

Extract more & less biased information
Data mining



Performance

Large/Deep NN

Conventional ML

Amount of Data

# Where we stand: Data revolution and cosmic evolution

**Understanding Cosmic evolution**

Data Surveys — Toolkit AI astro



recombination CMB

cosmic dawn CD 'first stars and galaxies'

reionisation EoR 'the last big phase transition'

cosmic noon peak of star formation

'late-time structure formation'

Years after the Big Bang

400 thousand · 0.1 billion · 1 billion · 4 billion · 8 billion · 13.8 billion

The Big Bang · Recombination · The Dark Age · Formation of first astronomical object

**JWST, HST**

**large galaxy surveys**

**Euclid, DESI**
**4MOST, Gaia**
Present day

**SKA-Low**

**SKA-Mid**

We are here!

D = 28 Gpc

'small galaxy surveys', wide-field line mapping

Fully ionized · Neutralized · Reionization · Fully ionized

z · ~1100 · 6 · 2 · 0

Credits: NAOJ, modified by C. Heneka

**Our goal:**
Learn about astrophysical & cosmological evolution
across cosmic time and scales

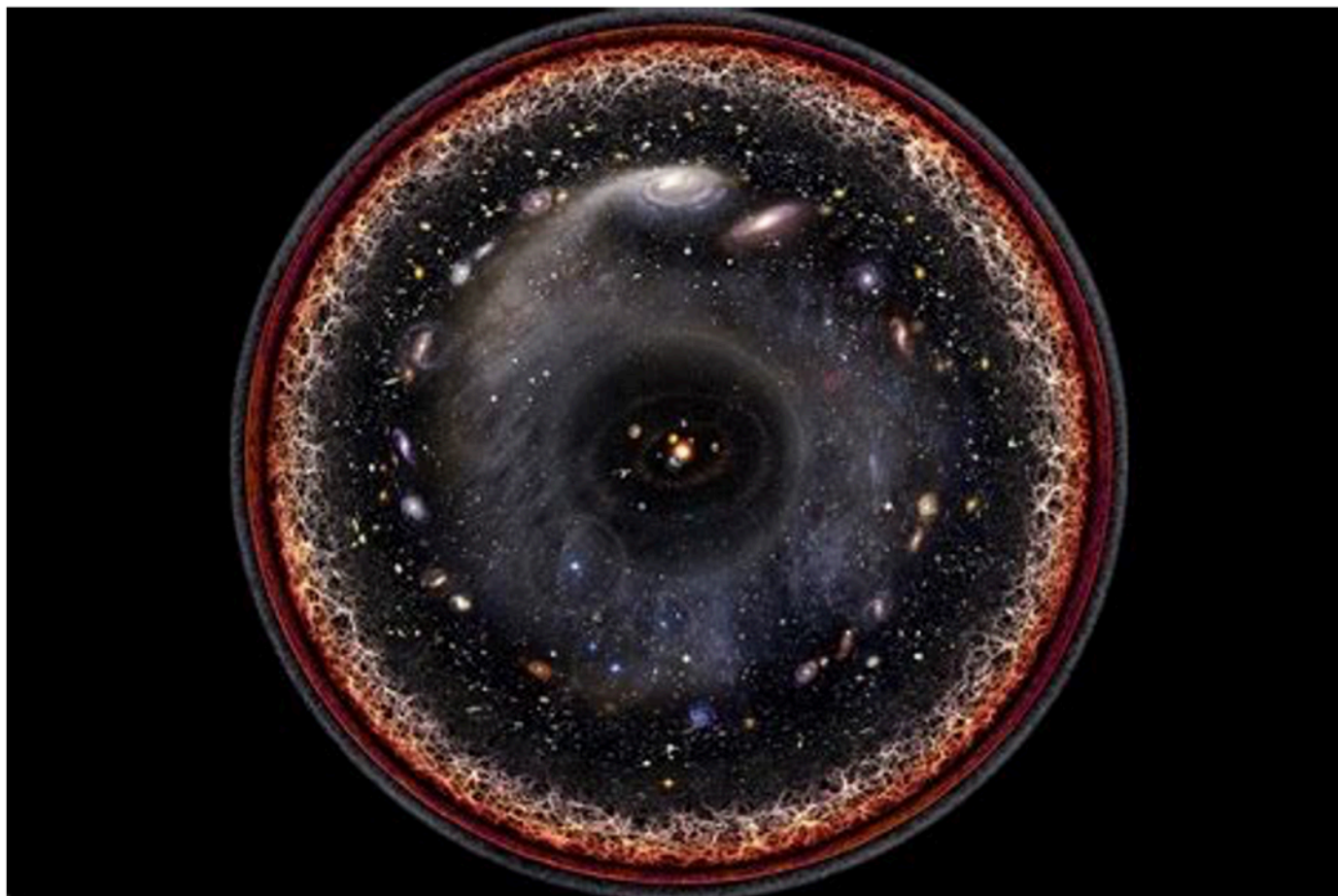Coming decade: push to map up to **80% of the observable Universe**

# … what does 80% of the observable Universe mean?

Modelling challenges

True LSS probes ⟶ orders of magnitude of scales up to the ultra-large

…what does 80% of the observable Universe even mean?



APOD, NASA, License & Credit: Wikipedia, Pablo Carlos Budassi

Observable Universe:
d ~ 28 Gpc (x3 Glyr)

80% if this:
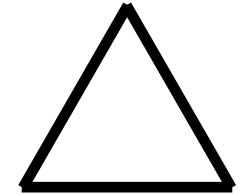d ~ 22 Gpc

Let's say we resolve (only) ~Mpc

⟶ about 3-4 orders of magnitude

⟶ about $10^9$-$10^{10}$ modes!

… at some point we sub-grid model and/or change modelling approach

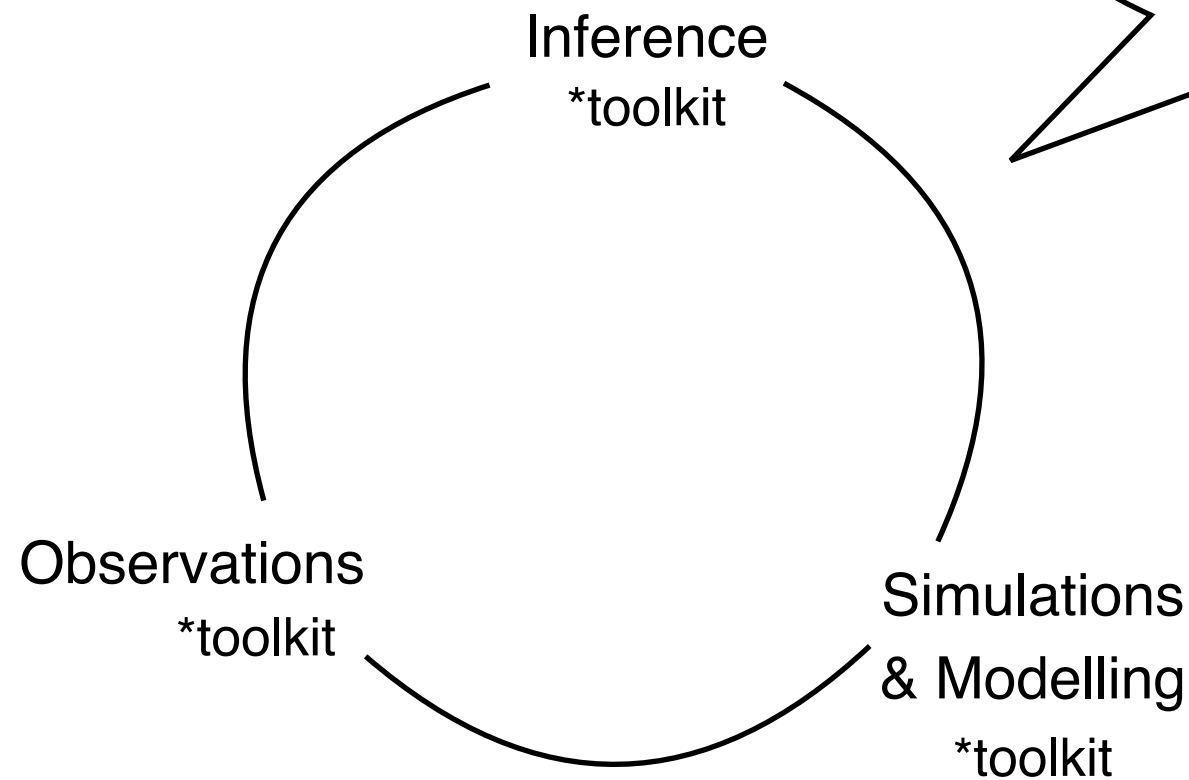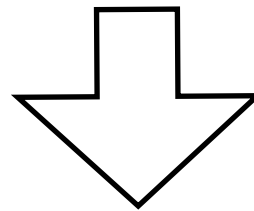How will Astrophysics and Cosmology advance in coming years?
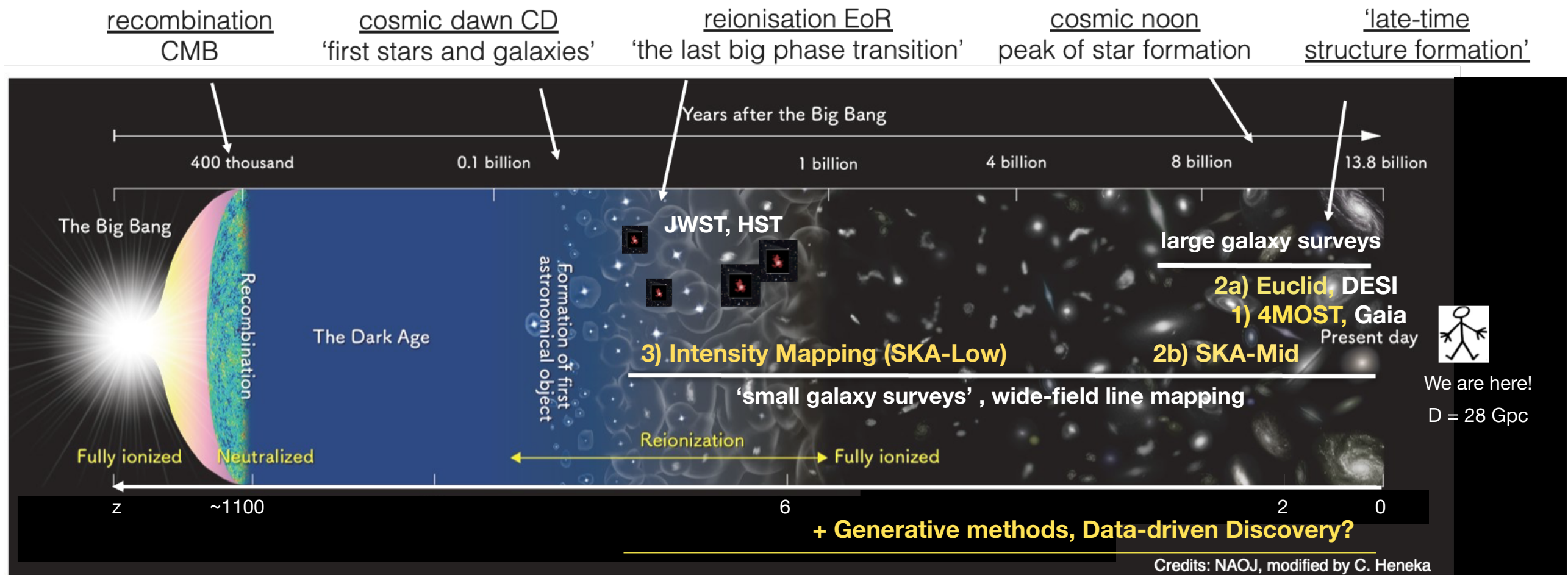
**Understanding Cosmic evolution**

**We need a versatile ML/AI toolkit*.**

New
Scientific Life Cycles

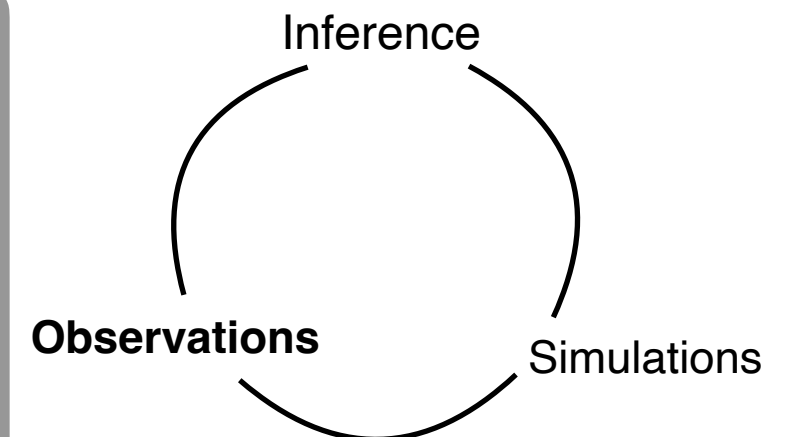Inference
*toolkit

**Examples**:
<u>Few sec:</u> Classification 40.000 spectra
<u>Few sec:</u> 7-parameter inference ~100MB cube
<u>Few sec:</u> detection, segmentation & flux
measurement O(100-1000) sources

Observations
*toolkit

Simulations
& Modelling
*toolkit

# Research Highlights: Astronomical Data Science and Artificial Intelligence

recombination
CMB

cosmic dawn CD
'first stars and galaxies'

reionisation EoR
'the last big phase transition'

cosmic noon
peak of star formation

'late-time
structure formation'

Years after the Big Bang

400 thousand — 0.1 billion — 1 billion — 4 billion — 8 billion — 13.8 billion

The Big Bang

JWST, HST

large galaxy surveys

2a) Euclid, DESI
1) 4MOST, Gaia
Present day

The Dark Age

3) Intensity Mapping (SKA-Low)     2b) SKA-Mid

We are here!
D = 28 Gpc

'small galaxy surveys' , wide-field line mapping

Reionization

Fully ionized   Neutralized

Fully ionized

z     ~1100     6     2     0

+ Generative methods, Data-driven Discovery?

Credits: NAOJ, modified by C. Heneka

## Select Highlights

1) **Classification**

2) Source detection & characterisation

3) Simulation-based inference
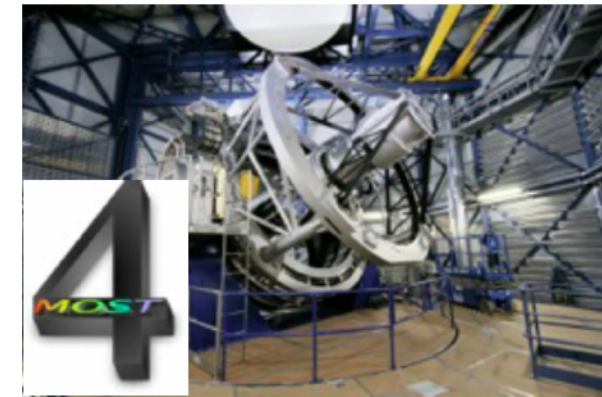
+ Generative methods, Data-driven Discovery

Inference

**Observations**

Simulations

# 1) Classification and triggering for large astronomical surveys

> **4MOST: On-the-fly classification of spectra (1D)**


https://www.4most.eu  Credit: ESO

- 5-year survey
- wide-field, fibre-fed, optical spectroscopy
- on ESO's 4-m-class telescope VISTA
- 2.5-degree diameter field-of-view, 2436 fibres
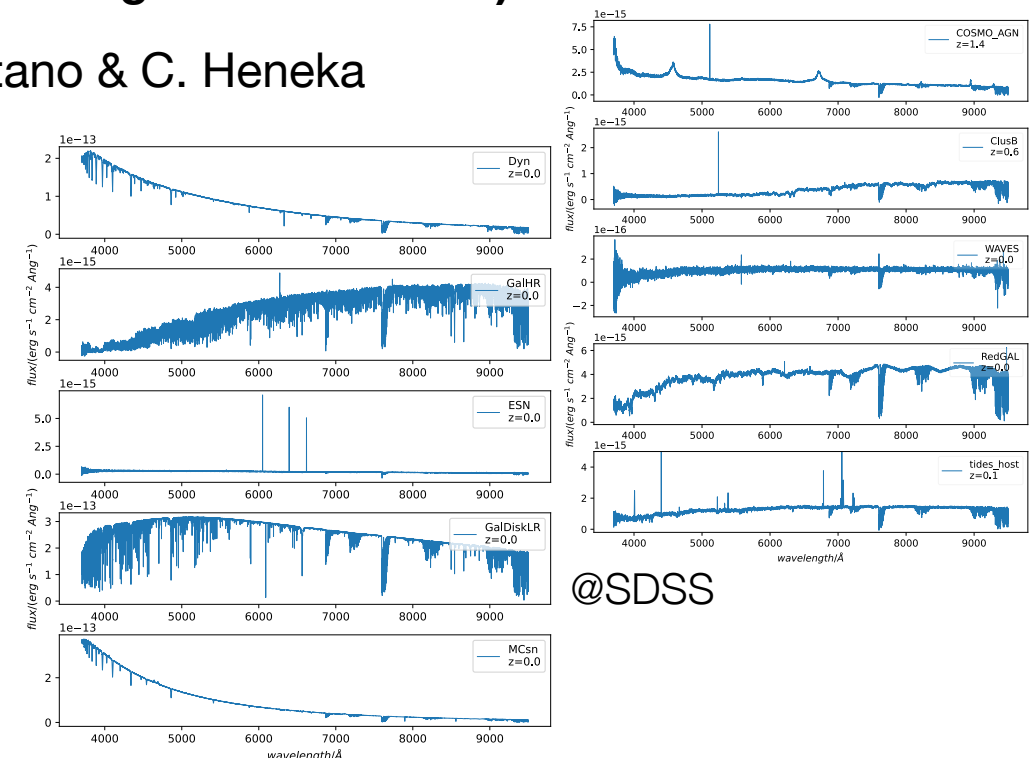- HRS R ≈ 18000 − 21000, LRS R ≈ 4000 − 7500
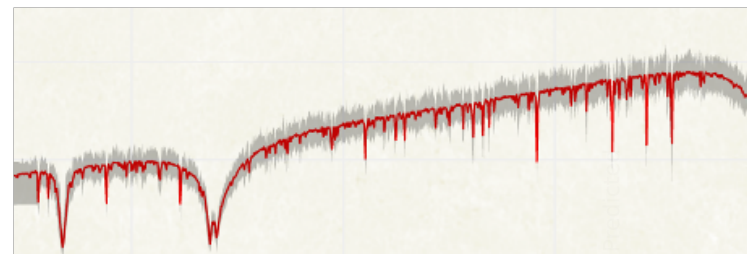- 20mio. (LRS), 3mio. (HRS) sources

**Goal: Data-driven classification pipeline layer (galactic & extragalactic sources)**

Classification infrastructure working group, led by: N. Napolitano & C. Heneka

⟶ *Benchmark with SDSS archival spectra:*


@SDSS

# 1) Classification and triggering for large astronomical surveys

**4MOST: On-the-fly classification of spectra (1D)**

**Goal: Data-driven classification pipeline layer (galactic & extragalactic sources)**

Classification infrastructure working group, led by: N. Napolitano & C. Heneka


https://www.4most.eu  Credit: ESO

Probabilistic multi-classifier

*For class:*
Convolutional
network variants

*For class uncertainties:*
Bayesian neural networks
and contrastive learning

*++ competitive with template fitting*

Examples:
**Few sec: Classification 40.000 spectra**
Few sec: 7-parameter inference ~100MB cube
Few sec: detection, segmentation & flux
measurement O(100-1000) sources

Zhong, Napolitano, Heneka+ arXiv:2311.04146

# Research Highlights: Astronomical Data Science and Artificial Intelligence



Credits: NAOJ, modified by C. Heneka

**Select Highlights**

1) Classification
2) **Source detection & characterisation**
3) Simulation-based inference
+ Generative methods, Data-driven Discovery

# 2a) The deblending problem

Example: Optical source detection & characterisation

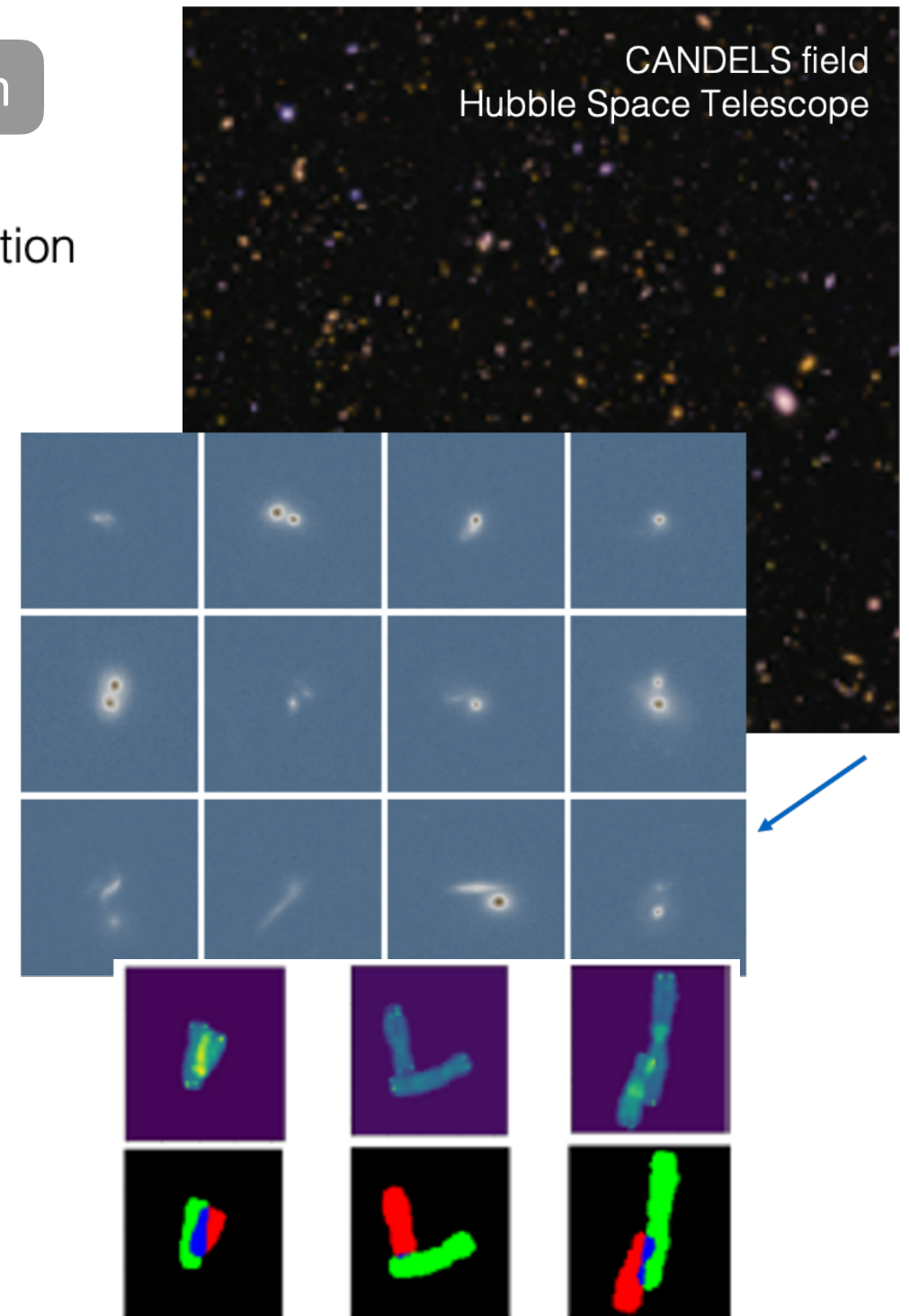**Goal:** 'Good' photometry for surveys with high blended fraction - avoid bias!

**Challenge:** Galaxies are 'transparent'

COIN network:
Emille Ishida (U. Clermont Auvergne)
Marc Huertas-Company (Obs. de Paris)
Alexandre Boucaud (APC, CNRS)

coindeblend

z1

z2



CANDELS field
Hubble Space Telescope

Boucaud, Huertas-Company, Heneka+ 20,
arXiv:1905.01324
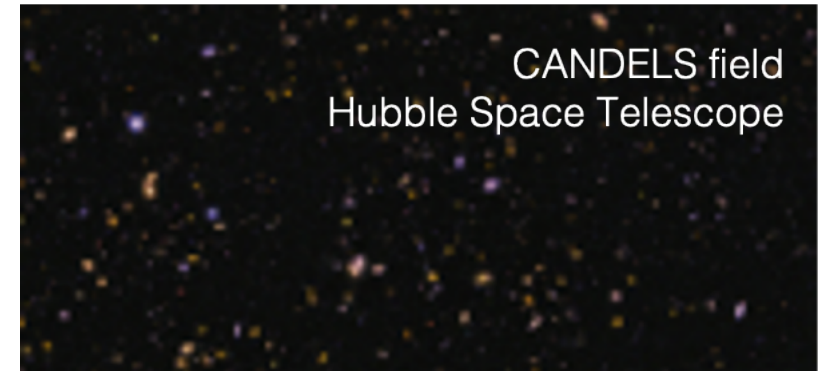
Lily Hu+ 2017

Similar challenge:
Overlapping chromosomes

# 2a) The deblending problem

Example: Optical source detection & characterisation

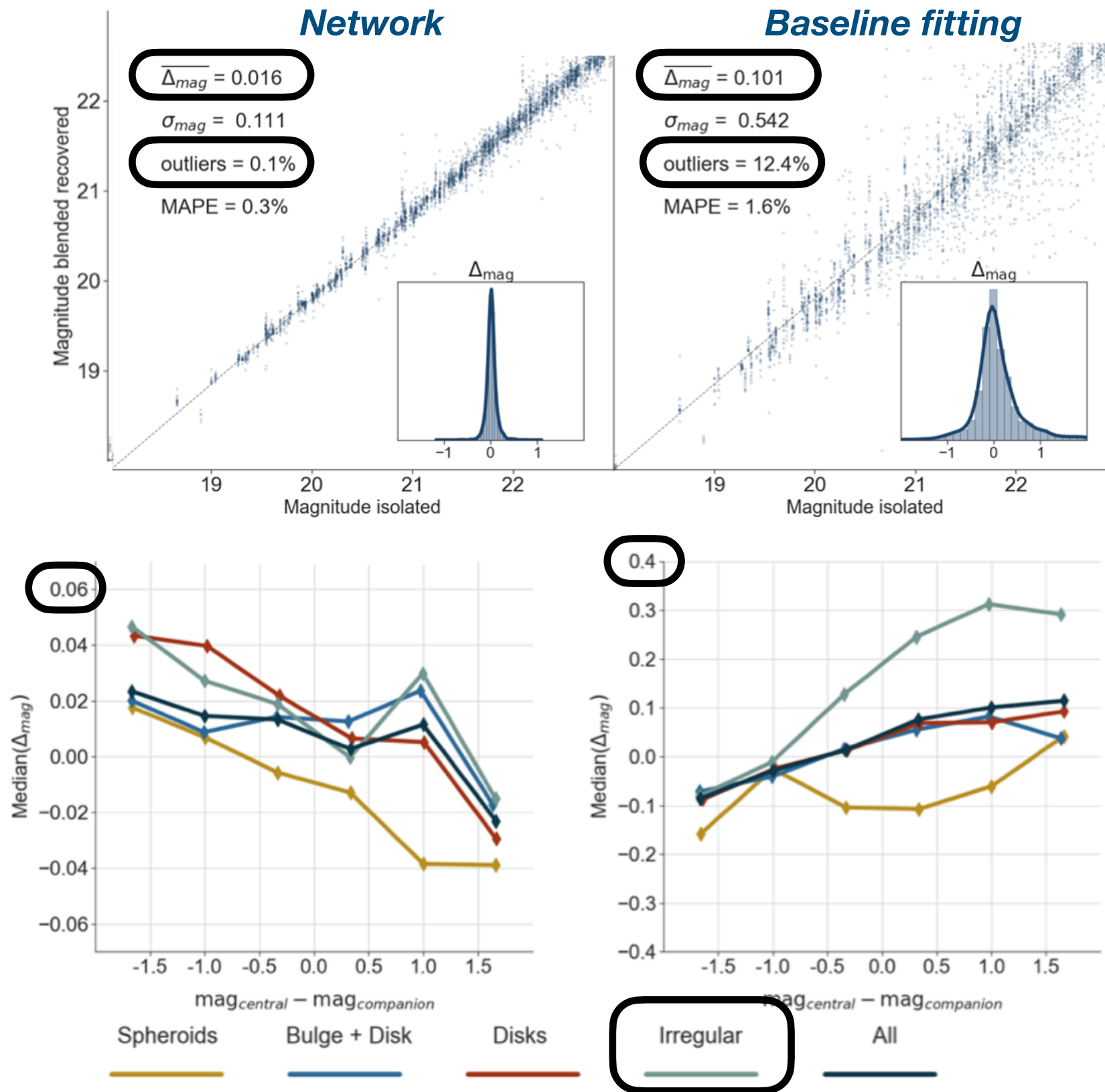**Goal:** 'Good' photometry for surveys with high blended fraction
- avoid bias!

Galaxy morphology

CANDELS field
Hubble Space Telescope

Credit: Euclid, ESA

# 2a) Optical source detection and characterisation



**Network**

$\overline{\Delta_{mag}} = 0.016$

$\sigma_{mag} = 0.111$

outliers = 0.1%

MAPE = 0.3%

**Baseline fitting**

$\overline{\Delta_{mag}} = 0.101$

$\sigma_{mag} = 0.542$

outliers = 12.4%

MAPE = 1.6%

Precise Photometry ✓

Bias ✓

Irregular ✓

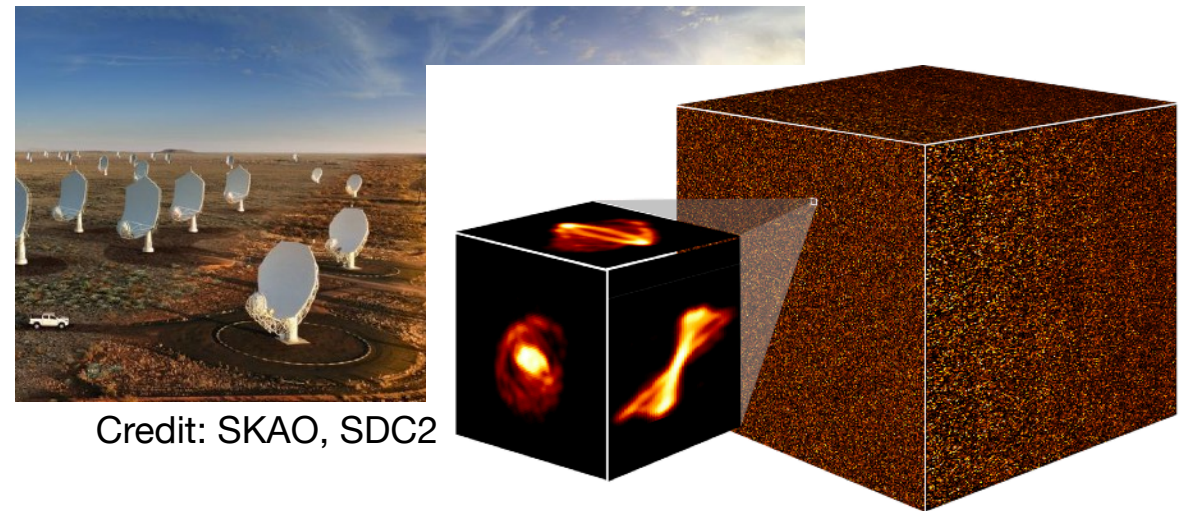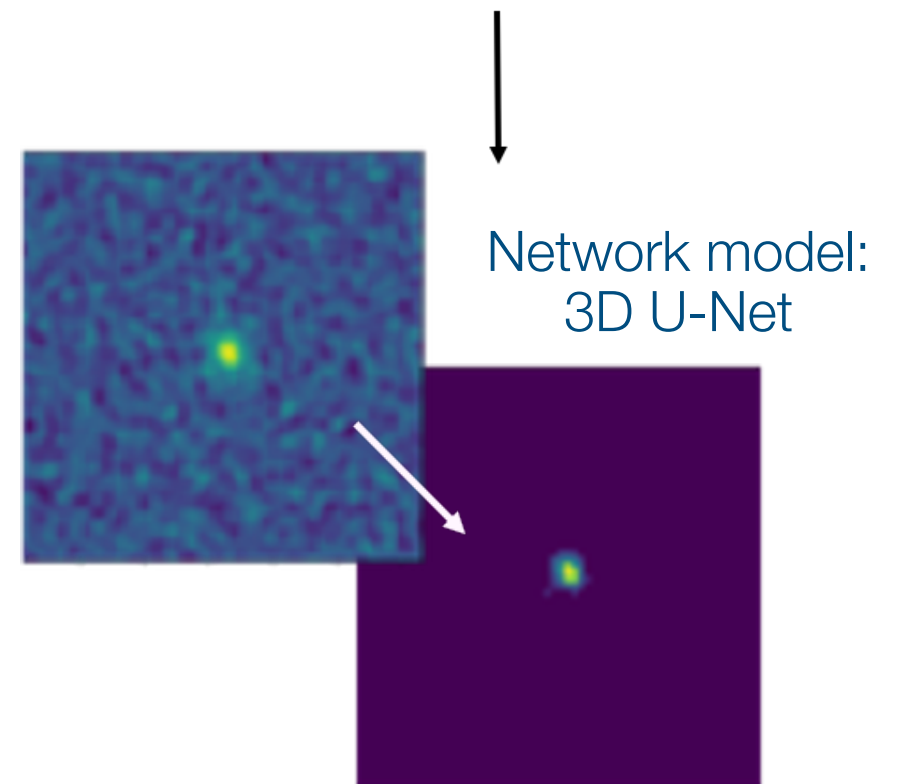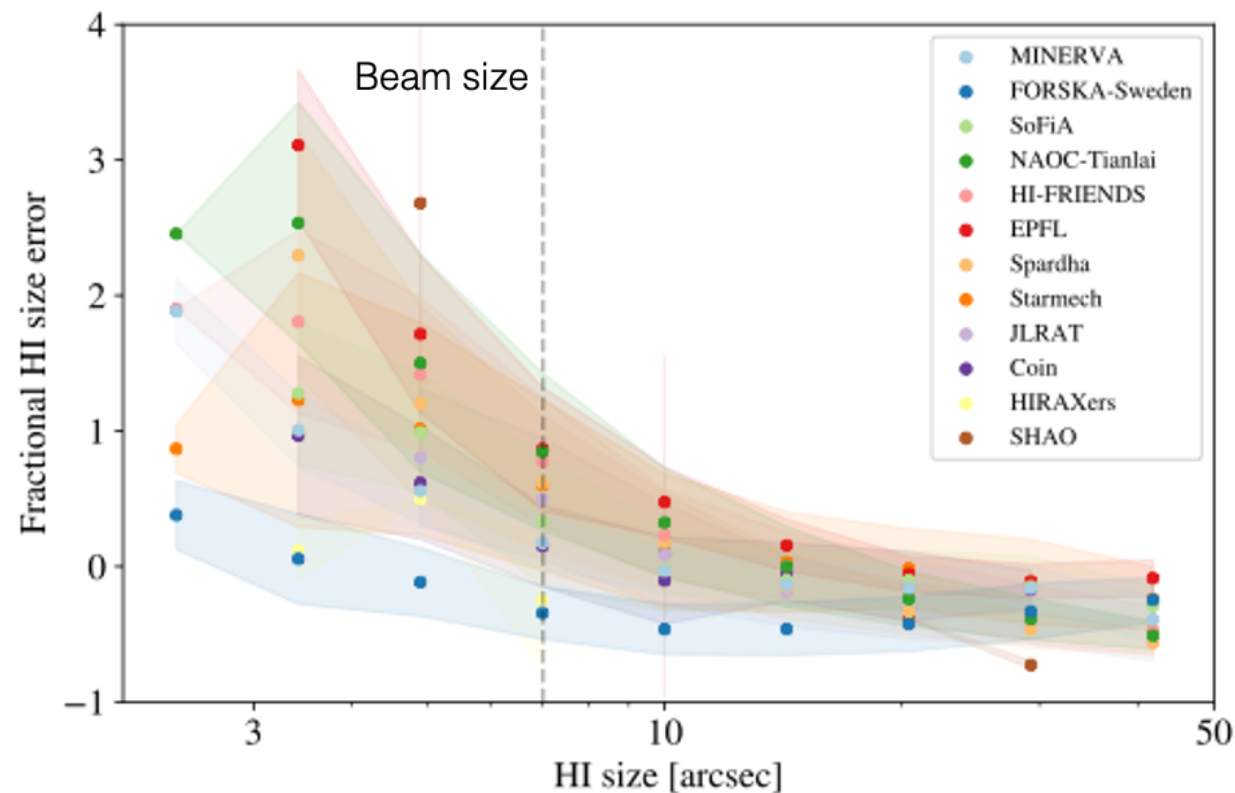coindeblend

# 2b) Radio source detection and characterisation

**Example: Source detection in tomographic data**

- 3D better than stitching of 2D + 1D
- High-fidelity 3D reconstructions
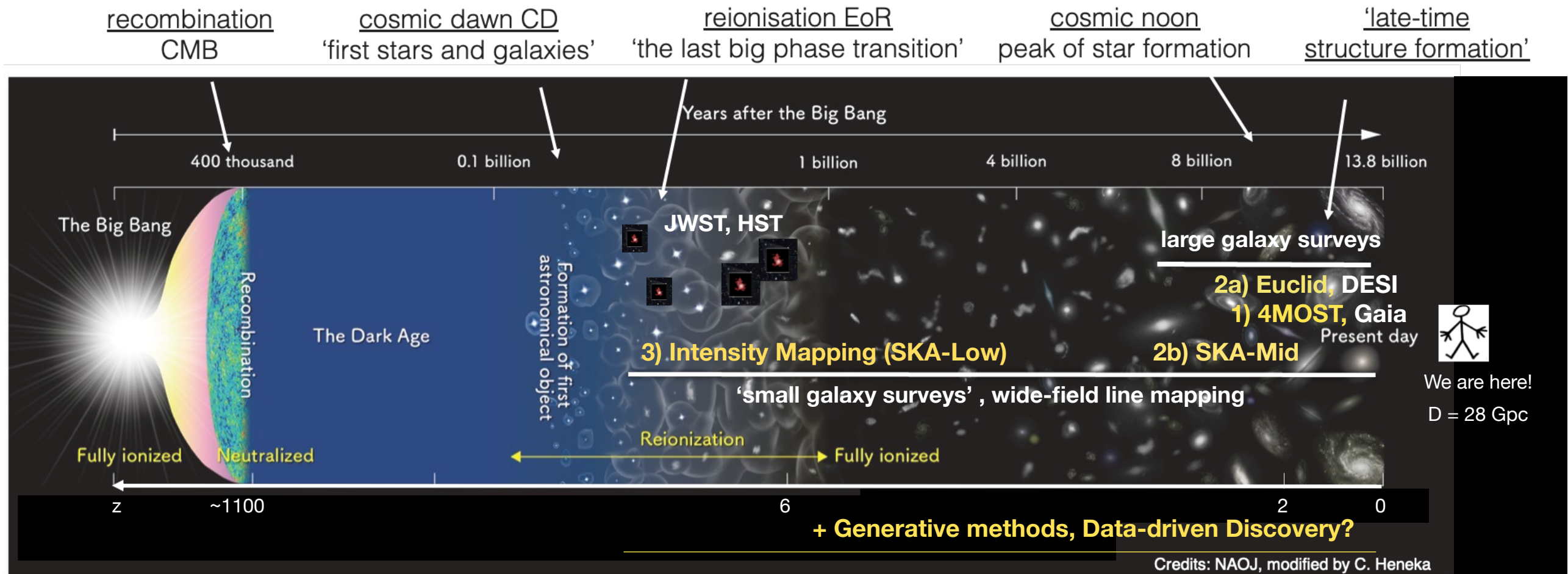- **Good prior for characterisation tasks via nets:**



Credit: SKAO, SDC2
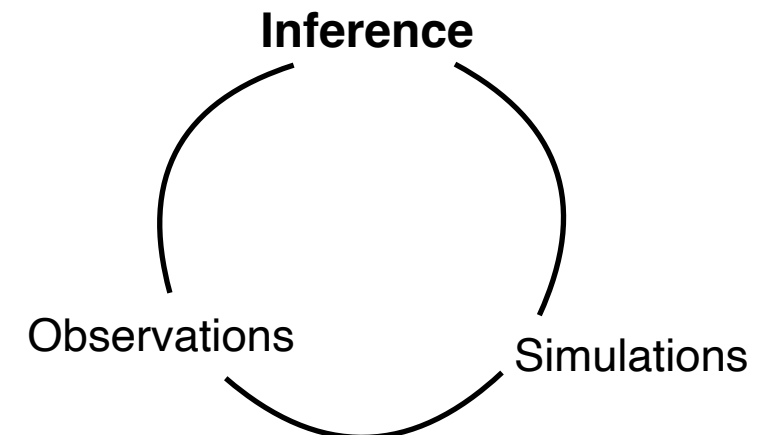
Total dimensions: (25,714 x 25,714 x 6,667) vox



Network model:
3D U-Net

@HPC/GPU Jean Zay (Idris)

Hartley+ 23 (incl. Heneka), arXiv:2303.07943
Heneka 23, arXiv:2311.17553

recombination
CMB

cosmic dawn CD
'first stars and galaxies'

reionisation EoR
'the last big phase transition'

cosmic noon
peak of star formation

'late-time
structure formation'

Years after the Big Bang

400 thousand          0.1 billion          1 billion          4 billion          8 billion          13.8 billion

The Big Bang

Recombination

The Dark Age

Formation of first astronomical object

JWST, HST

large galaxy surveys

2a) Euclid, DESI
1) 4MOST, Gaia
Present day

3) Intensity Mapping (SKA-Low)          2b) SKA-Mid

'small galaxy surveys' , wide-field line mapping

We are here!

D = 28 Gpc

Fully ionized   Neutralized

Reionization

Fully ionized

z        ~1100                          6                                    2        0

+ Generative methods, Data-driven Discovery?

Credits: NAOJ, modified by C. Heneka

## Select Highlights

1) Classification
2) Source detection & characterisation
3) **Simulation-based inference**
+ Generative methods, Data-driven Discovery

**Inference**

Observations                Simulations

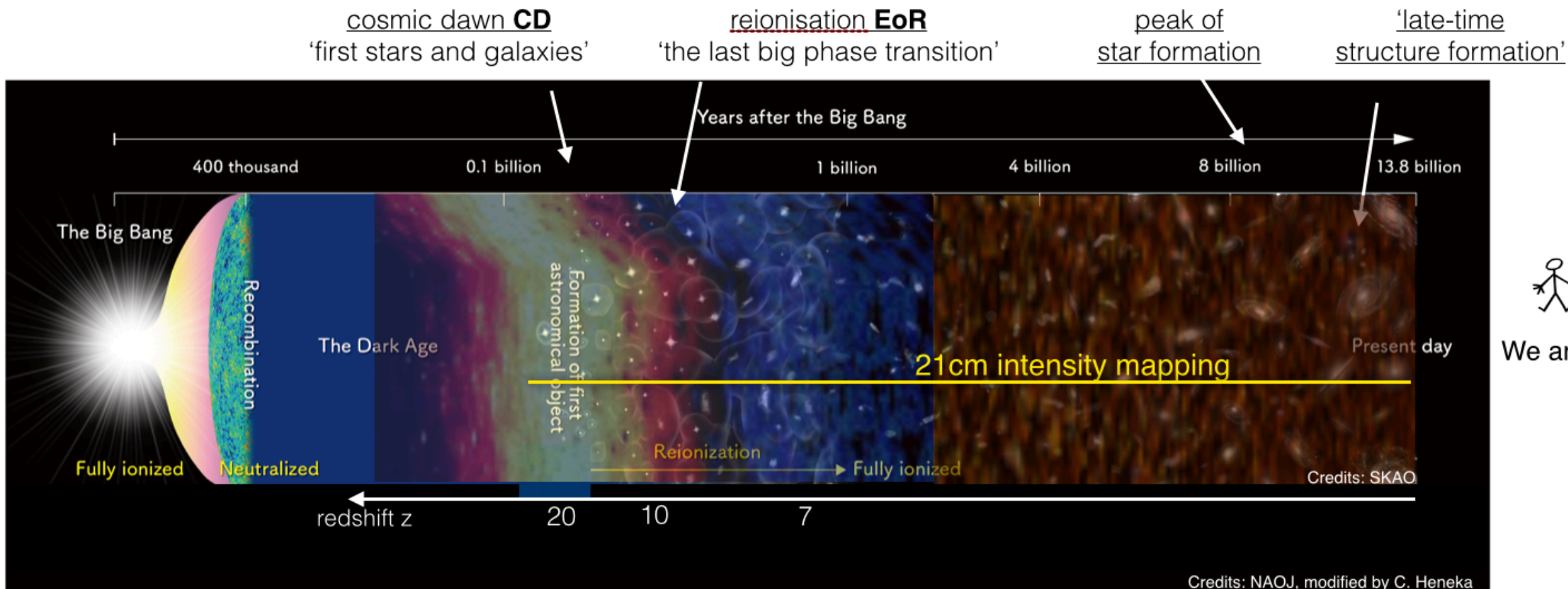# 3) 3D Simulation-based inference (SBI) for the SKA



@SKAO

SKA: TB/s, few EB/day
Archive: ~700 PB/yr

Why care?

Tomography of >50% of the Universe

True 'Big Data'
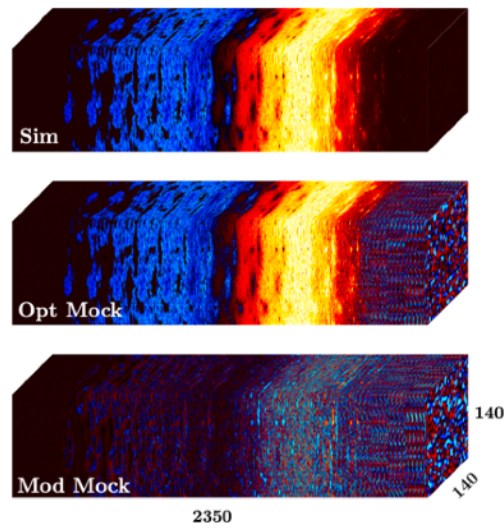
non-linear, non-Gaussian signal

⟶ MCMC becomes slow and biased
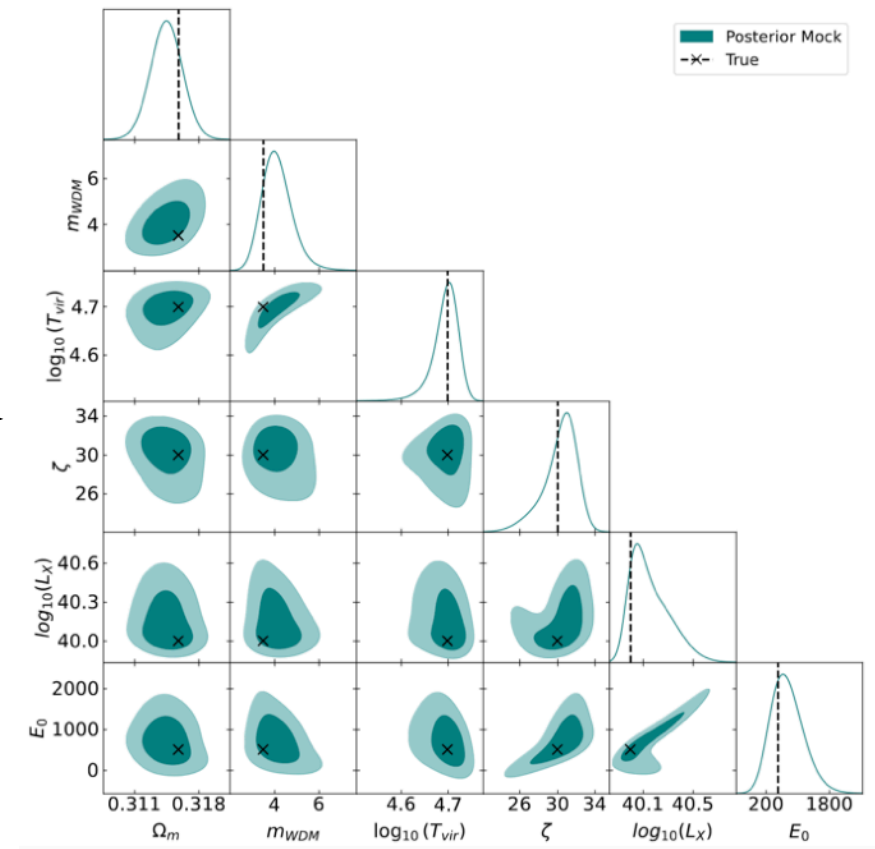
⟶ Move to full likelihood(-free) inference with networks
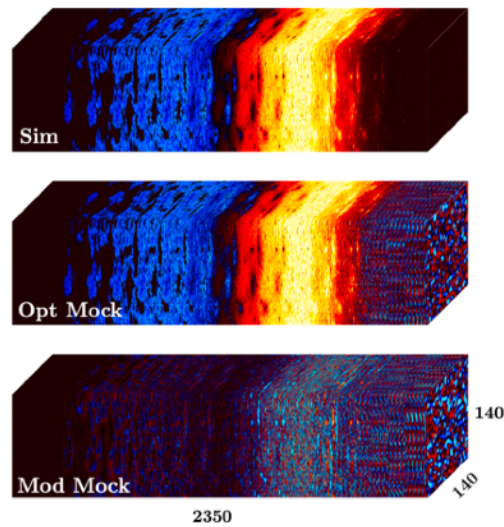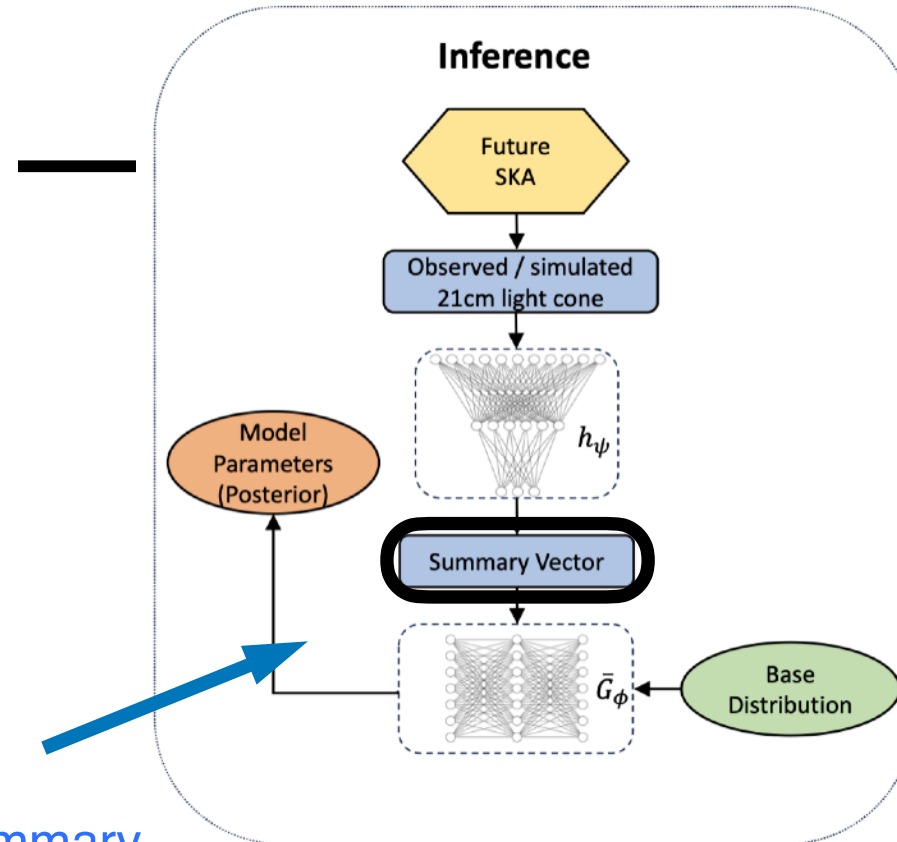
# 3) Simulation-based inference (SBI) for the SKA



Neutsch, Heneka, Brüggen (2022), arXiv:2201.07587

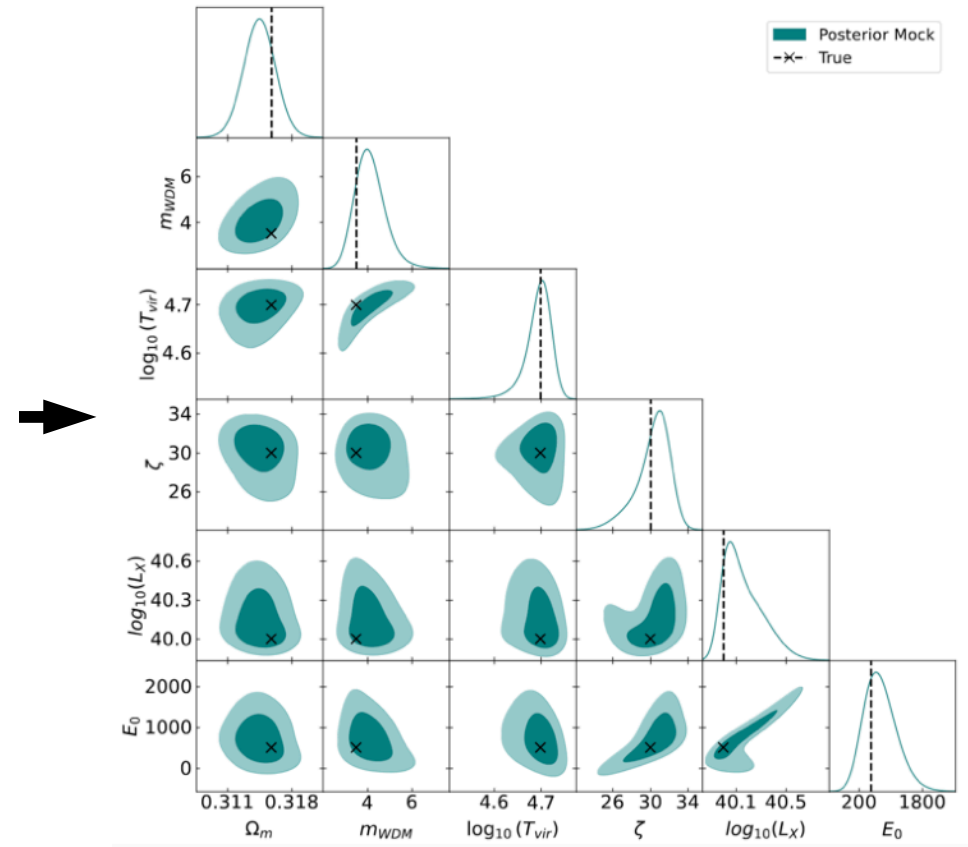Schosser, Heneka, Plehn, arXiv:2401.04174

# 3) Simulation-based inference (SBI) for the SKA



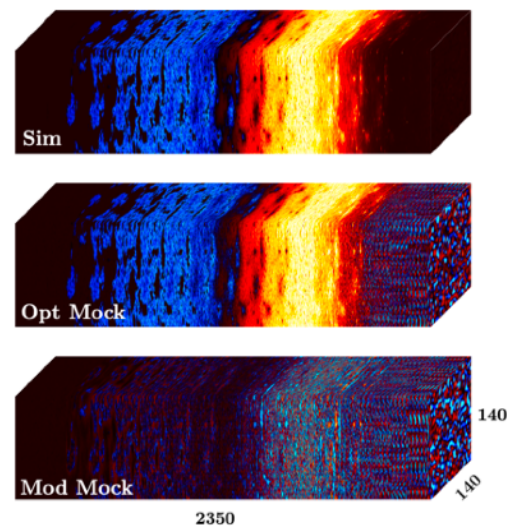**SBI with flows/cINN**

based on BayesFlow

arXiv:2003.06281

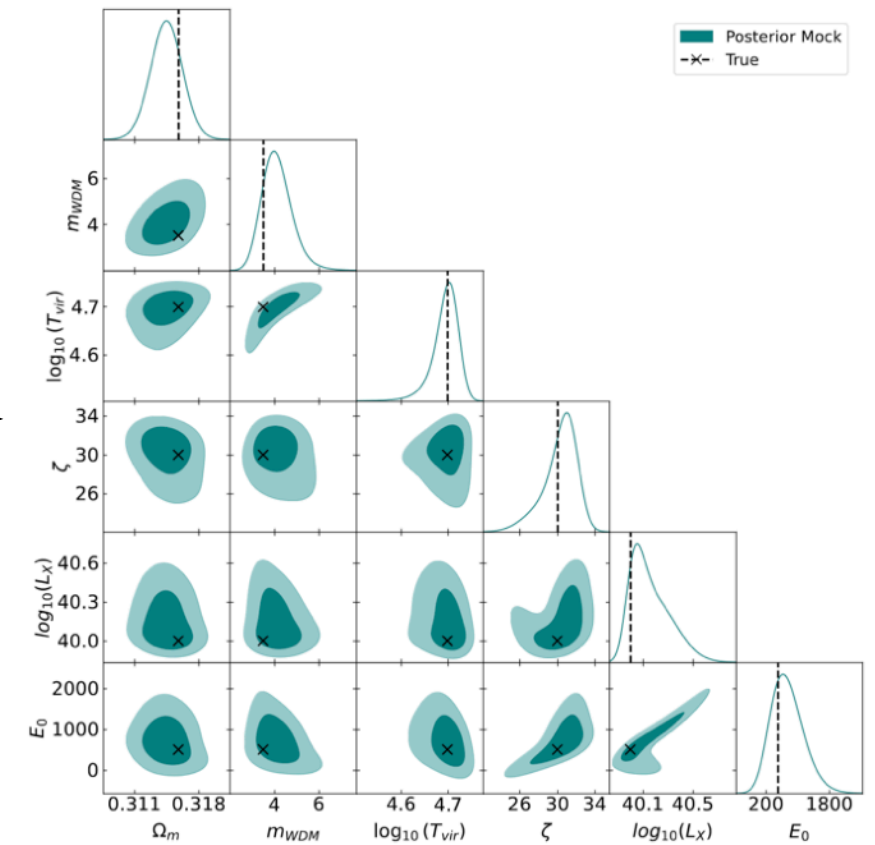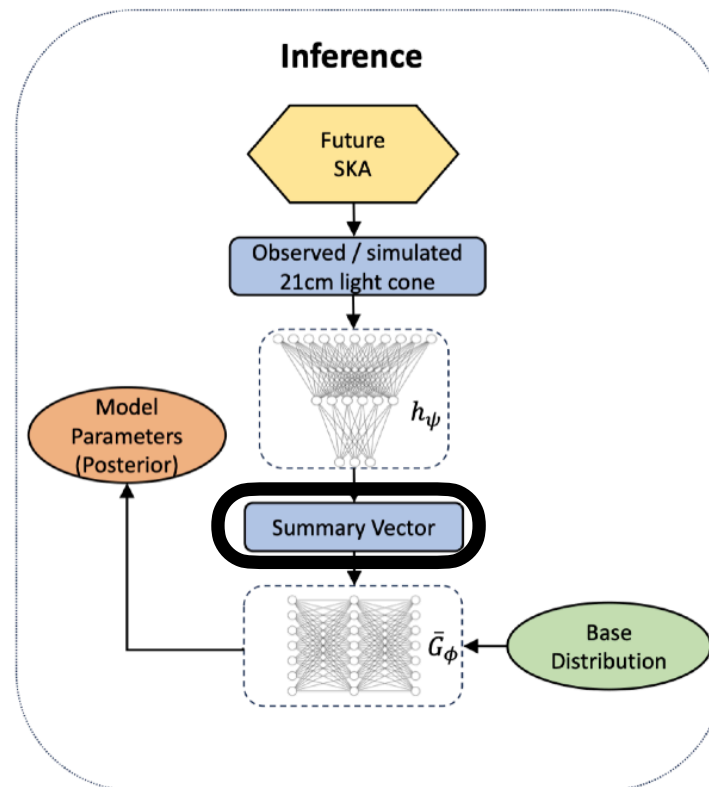'The cool' here:

End-to-end training

Learned optimal summary

Neutsch, Heneka, Brüggen (2022), arXiv:2201.07587

Schosser, Heneka, Plehn, arXiv:2401.04174

# 3) Simulation-based inference (SBI) for the SKA



SBI with flows/cINN

based on BayesFlow

arXiv:2003.06281

Sim: Summary stays close to original

Mock: Heavy adjustment of summary vector

We profit from learned summary in presence of noise (more)!

Neutsch, Heneka, Brüggen (2022), arXiv:2201.07587
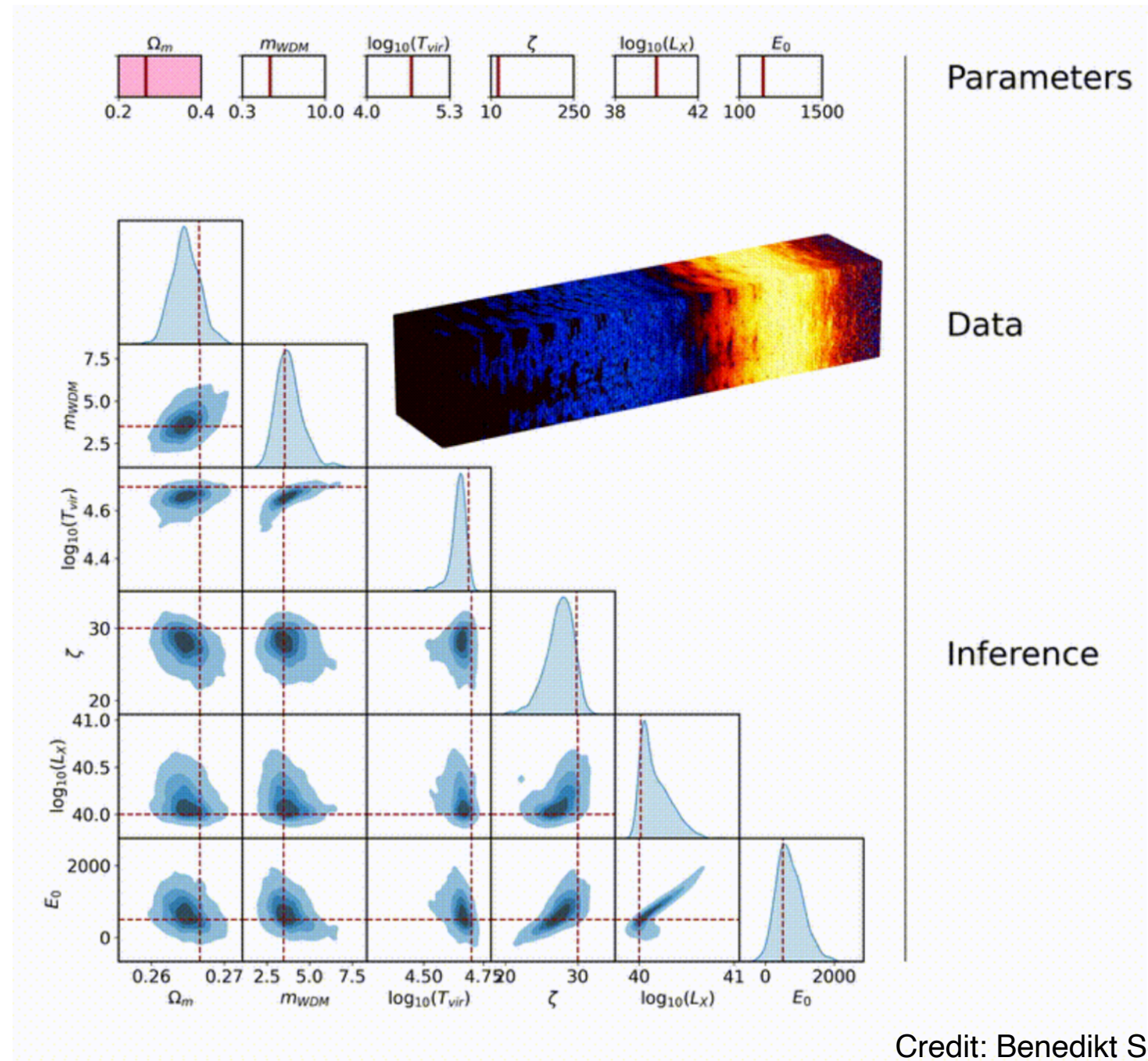
Schosser, Heneka, Plehn, arXiv:2401.04174

# 3) Simulation-based inference (SBI) for the SKA

Performance validation via:

- Distribution of latent variables
- Simulation-based calibration
- Parameter recovery
- Mutual information

**Trained SBI in action:**

**1 frame = 1 MCMC**

astro-ML

21cm_pie  (Public)



Credit: Benedikt Schosser

**'Optimal, fast, and robust inference of reionization-era cosmology with the 21cmPIE-INN'**

Schosser, Heneka, Plehn (2024), arXiv:2401.04174

Performance validation via:

- Distribution of latent variables

- Simulation-based calibration

- Parameter recovery

- Mutual information

**Trained SBI in action:**

**1 frame = 1 MCMC**

astro-ML

21cm_pie (Public)



Credit: Benedikt Schosser

**'Optimal, fast, and robust inference of reionization-era cosmology with the 21cmPIE-INN'**

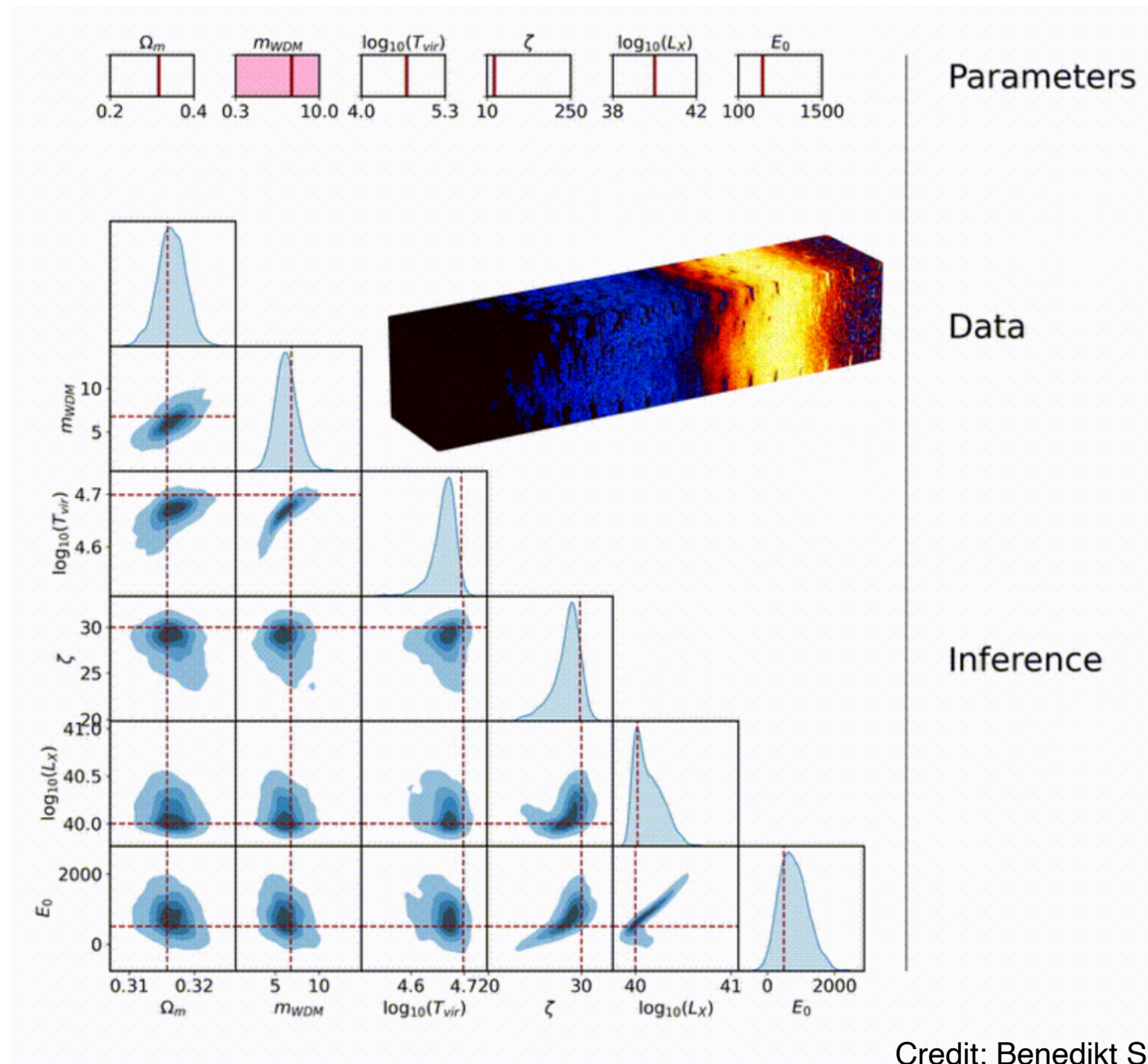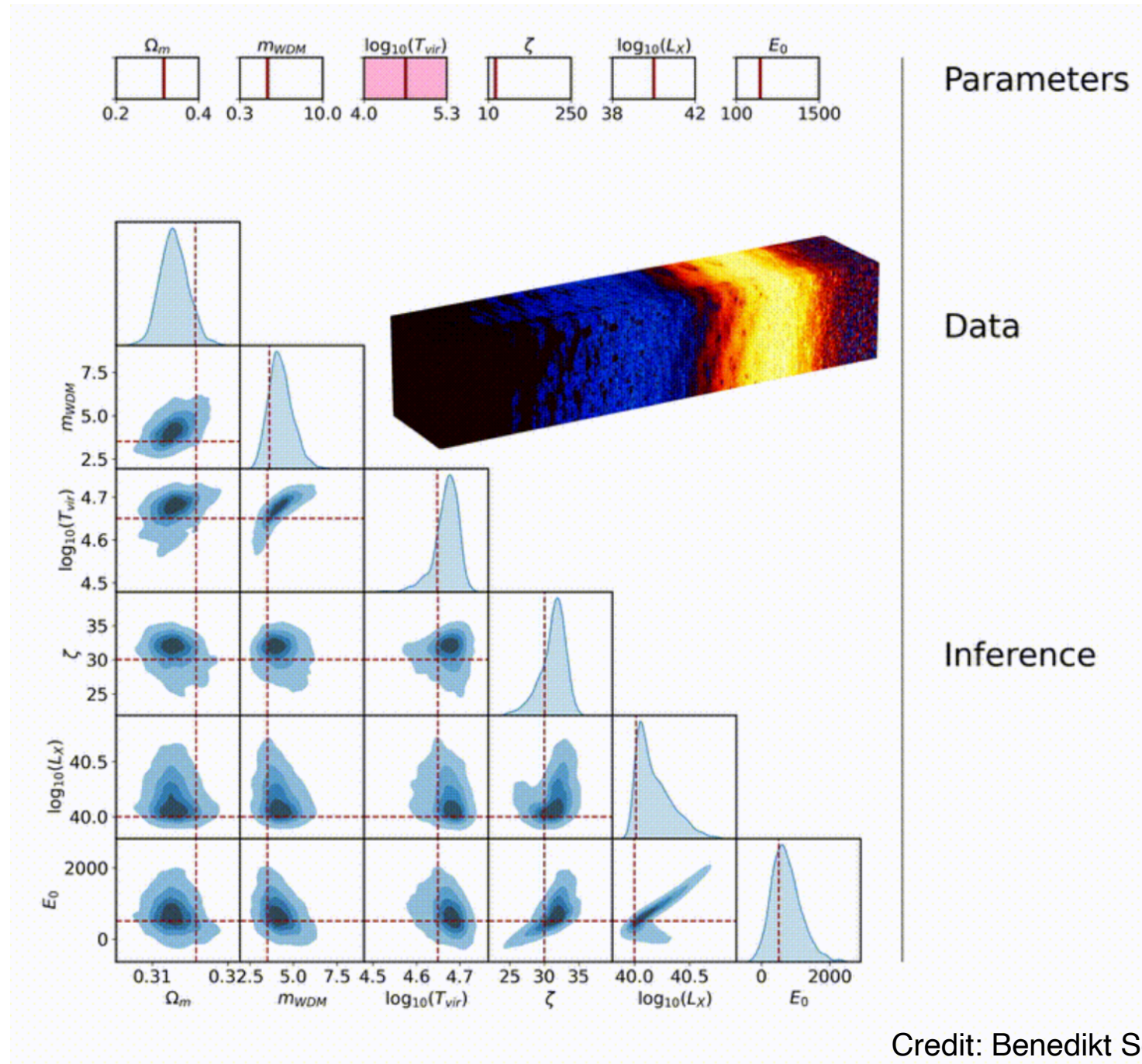Schosser, Heneka, Plehn (2024), arXiv:2401.04174

# 3) Simulation-based inference (SBI) for the SKA

Performance validation via:

- Distribution of latent variables
- Simulation-based calibration
- Parameter recovery
- Mutual information

**Trained SBI in action:**

**1 frame = 1 MCMC**

astro-ML

21cm_pie  (Public)



Credit: Benedikt Schosser

**'Optimal, fast, and robust inference of reionization-era cosmology with the 21cmPIE-INN'**
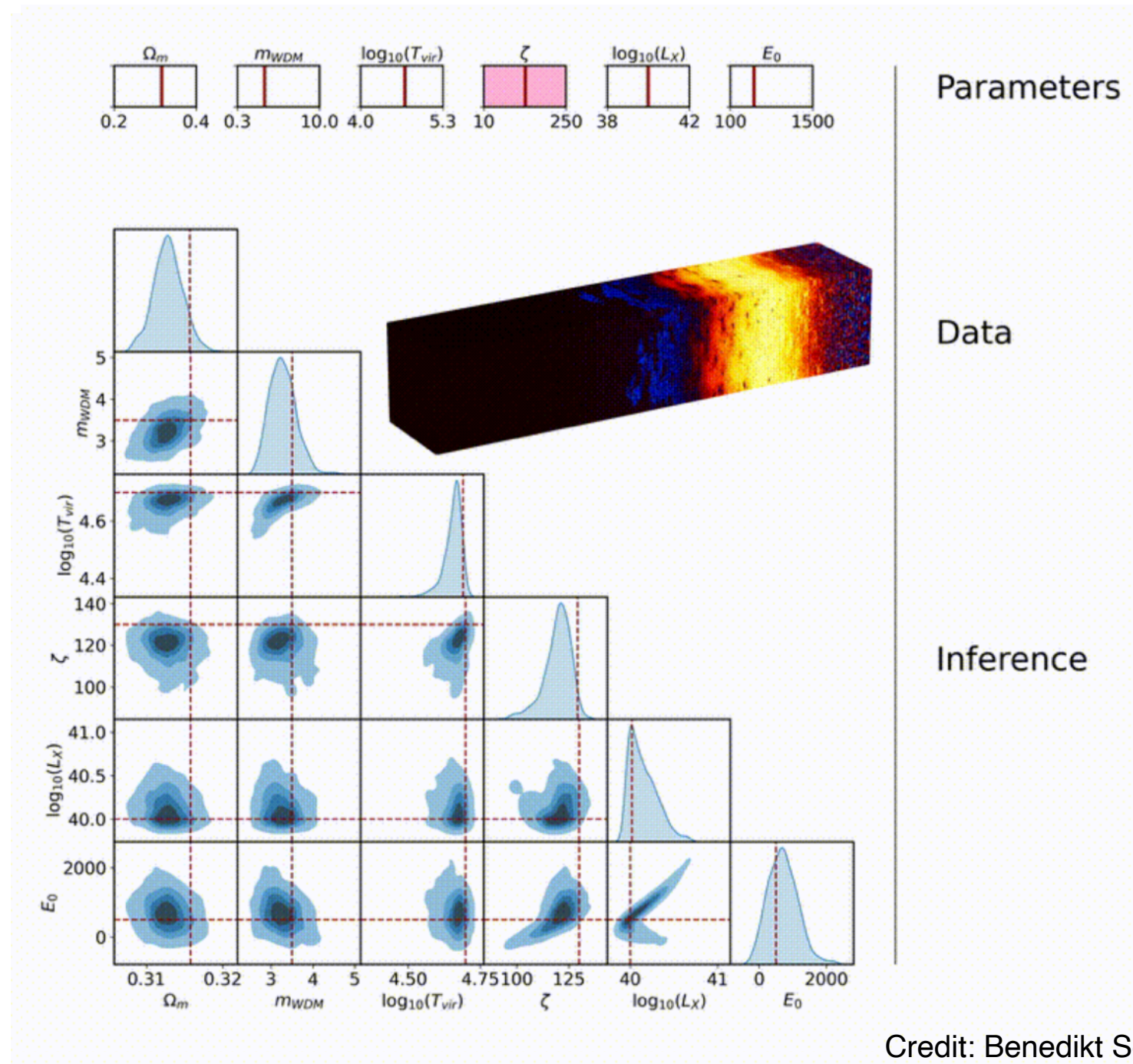
Schosser, Heneka, Plehn (2024), arXiv:2401.04174

# 3) Simulation-based inference (SBI) for the SKA

Performance validation via:

- Distribution of latent variables
- Simulation-based calibration
- Parameter recovery
- Mutual information

**Trained SBI in action:**

**1 frame = 1 MCMC**

astro-ML

21cm_pie (Public)



Credit: Benedikt Schosser

**'Optimal, fast, and robust inference of reionization-era cosmology with the 21cmPIE-INN'**

Schosser, Heneka, Plehn (2024), arXiv:2401.04174

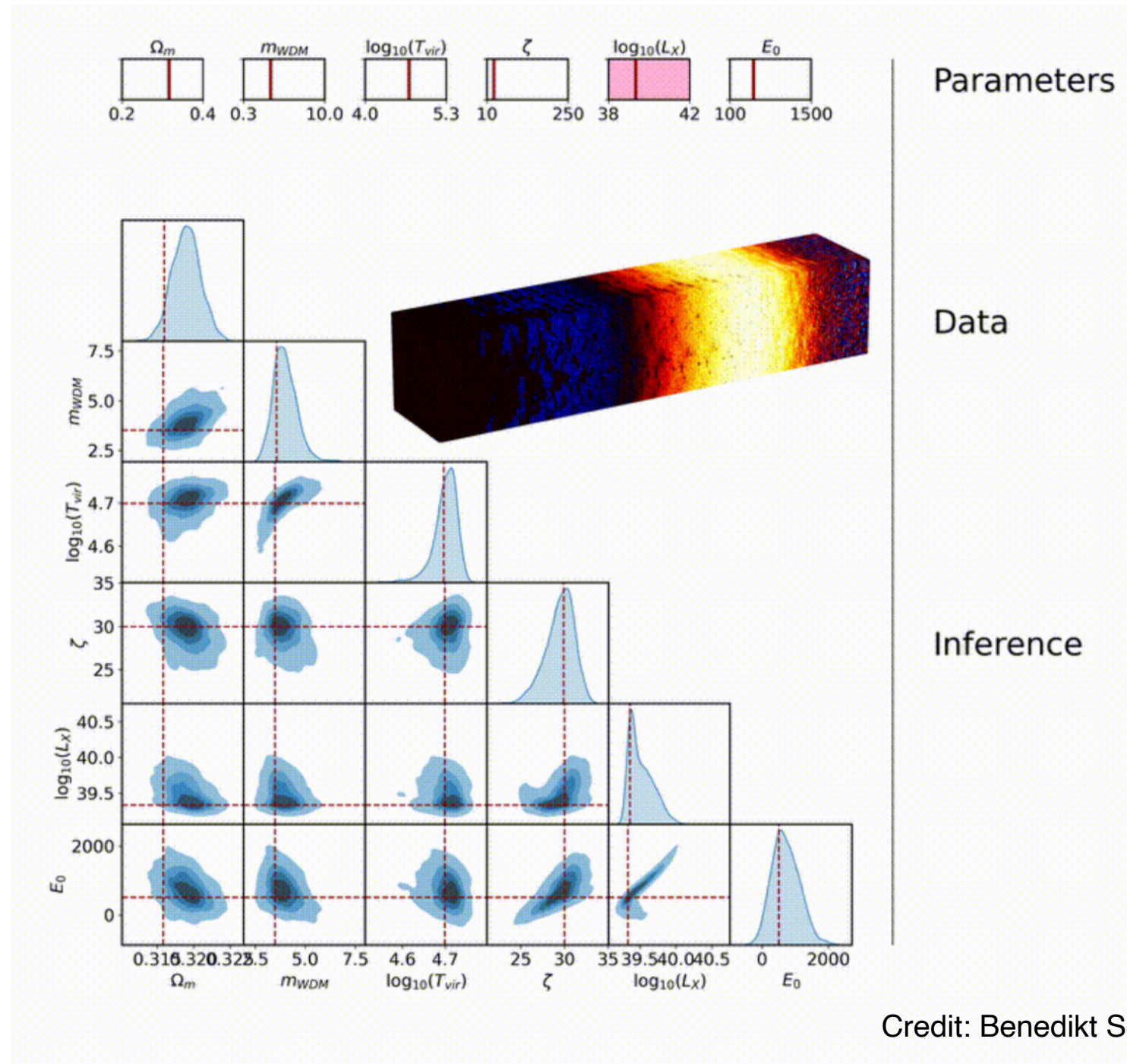# 3) Simulation-based inference (SBI) for the SKA

Performance validation via:

- Distribution of latent variables
- Simulation-based calibration
- Parameter recovery
- Mutual information

**Trained SBI in action:**

**1 frame = 1 MCMC**

astro-ML

🟩 21cm_pie  Public



Credit: Benedikt Schosser

**'Optimal, fast, and robust inference of reionization-era cosmology with the 21cmPIE-INN'**
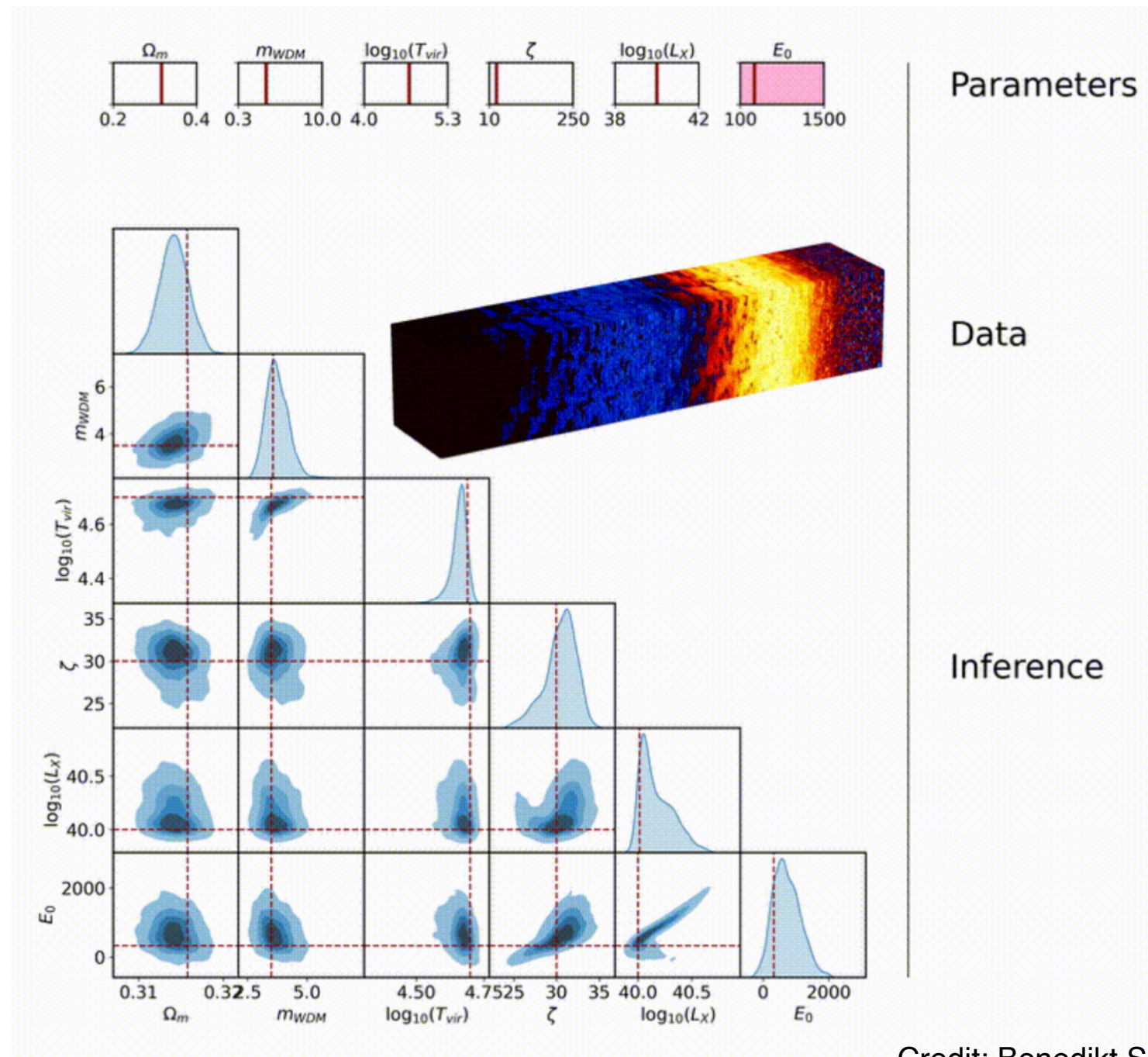
Schosser, Heneka, Plehn (2024), arXiv:2401.04174

# 3) Simulation-based inference (SBI) for the SKA

Performance validation via:

- Distribution of latent variables
- Simulation-based calibration
- Parameter recovery
- Mutual information

**Trained SBI in action:**

**1 frame = 1 MCMC**
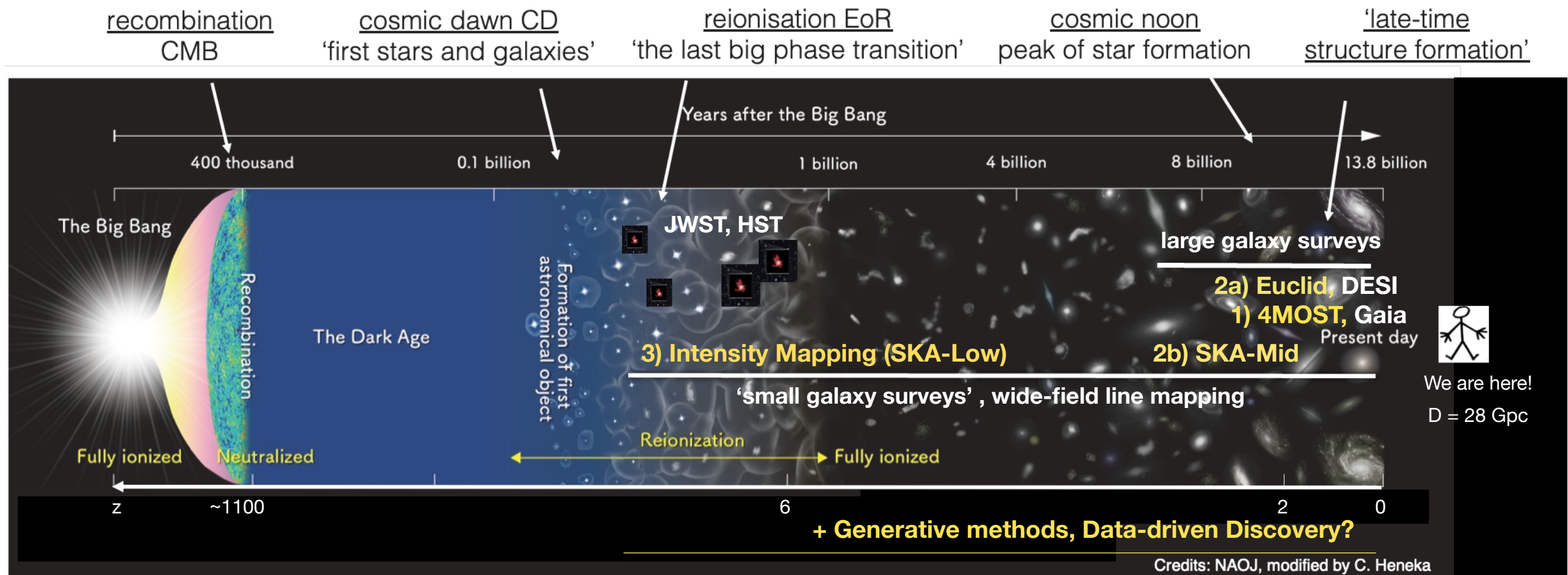
astro-ML

21cm_pie (Public)



Credit: Benedikt Schosser

**'Optimal, fast, and robust inference of reionization-era cosmology with the 21cmPIE-INN'**

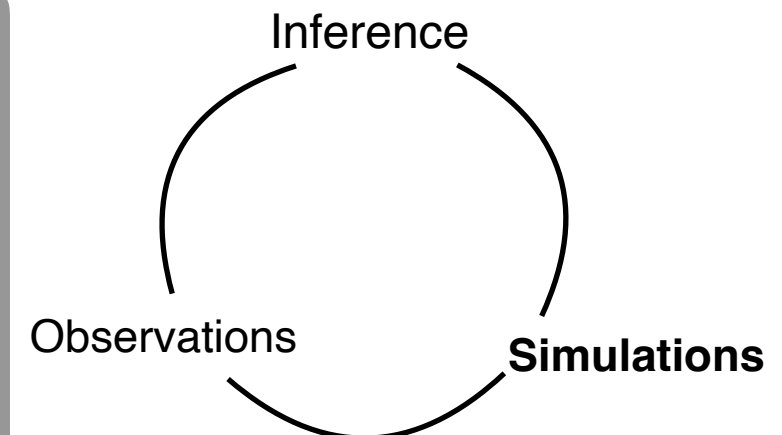Schosser, Heneka, Plehn (2024), arXiv:2401.04174

Credits: NAOJ, modified by C. Heneka

**Select Highlights**

1) Classification
2) Source detection & characterisation
3) Simulation-based inference
+ **Generative methods, Data-driven Discovery**

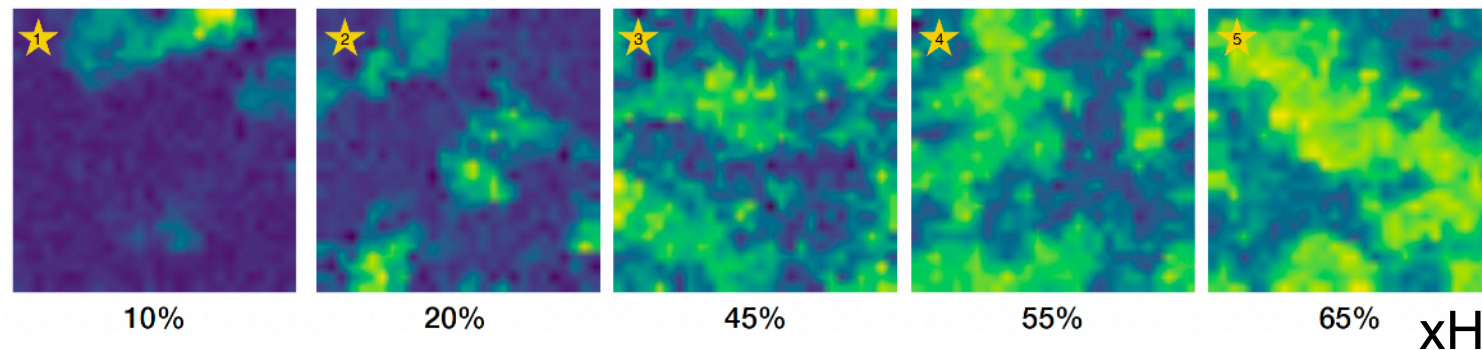# + Generative Modelling, Data-driven Discovery

Generative models:

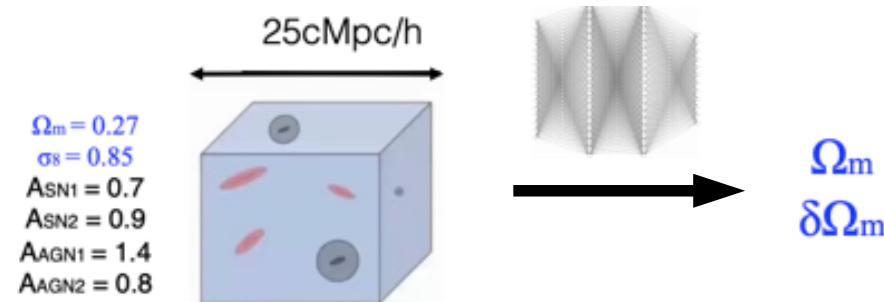Is there a fast way to emulate whole simulations?



| 10% | 20% | 45% | 55% | 65% xH |

**Diffusion models** (Ho+20)

Trained on 21cmFAST
(Mesinger+12, Murray+20)

@Lara Alegre (Postdoc ITP)

+



25cMpc/h

$\Omega_m = 0.27$
$\sigma_8 = 0.85$
$A_{SN1} = 0.7$
$A_{SN2} = 0.9$
$A_{AGN1} = 1.4$
$A_{AGN2} = 0.8$

$\Omega_m$
$\delta\Omega_m$

YES!

~10% uncertainty

Connections on a high-dim manifold?

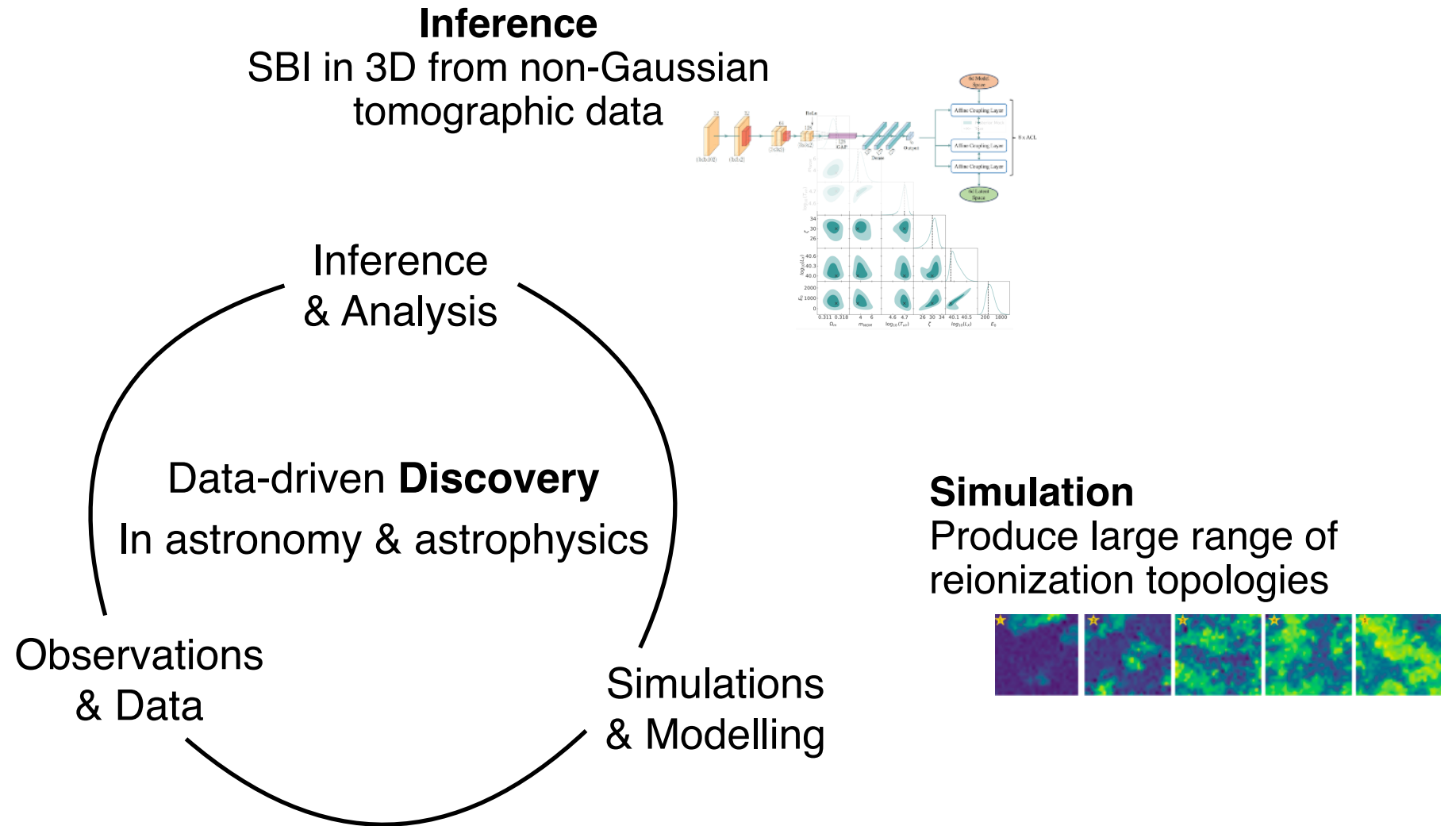Villaescusa-Navarro (incl. Heneka)+22

**Data-driven discovery:**

Can we measure $\Omega_m$ only from one (random) galaxy?

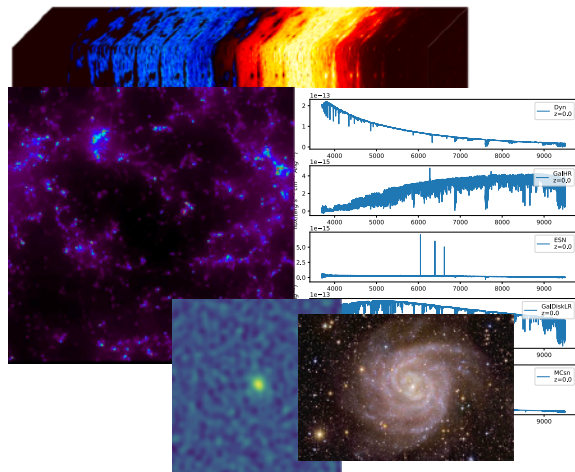… on the way to scientific discovery with ML/AI/Big Data and the SKAO !?

# Summary: Where we stand

**Goal: Understanding & discovery**

**Inference**
SBI in 3D from non-Gaussian
tomographic data

Inference
& Analysis

Data-driven **Discovery**
In astronomy & astrophysics

Observations
& Data

Simulations
& Modelling

**Detection & Characterisation**
Unbiased measurements from
diverse sources (galaxies)

**Classification**
Online classifier and triggering

**Simulation**
Produce large range of
reionization topologies

Select publications:
Neutsch, Heneka, Brüggen 22, arXiv:2201.07587
Schosser, Heneka, Plehn, arXiv:2401.04174
Hartley+ 23 (incl. Heneka), arXiv:2303.07943
Heneka 23, arXiv:2311.17553
Boucaud, Huertas-Company, Heneka+ 20, arXiv:1905.01324
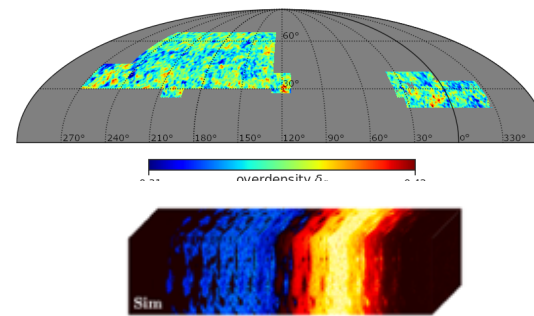Zhong, Napolitano, Heneka+24,  arXiv:2311.04146

**Next stop**: Robust Foundational models



Generation & modelling

Data-Simulation Gap

1) **Data, Mocks**

Interferometric observations

Mock observations

Maps / tomography

2) **Transfer**

**Simulation-based:**

Optimal Summaries

Foundational ?

tuning, self-supervised importance weighing

Stay tuned: Ayo Ore et al.

3) **Inference**

Currently:

Comparison to 'random mocks'
Derive summaries, such as $C_\ell^{gg}$

\+ jackknife

\+ MCMC sampling

# Summary and conclusions

ML/DL Physics publications on arXiv



Some advertisement:

We have 'endless' open-source public data

Astronomy & Astrophysics as ideal 'playground' for DL-driven discovery

Percent Astrophsics (Physics) ML/DL publications on arXiv



*Thank you for your attention!*
*heneka@thphys.uni-heidelberg.de*